

# Vision-to-Language Generation

---

XIN WANG  
UCSB CS281B

# What is Visual Recognition?

---



# What is Visual Recognition?

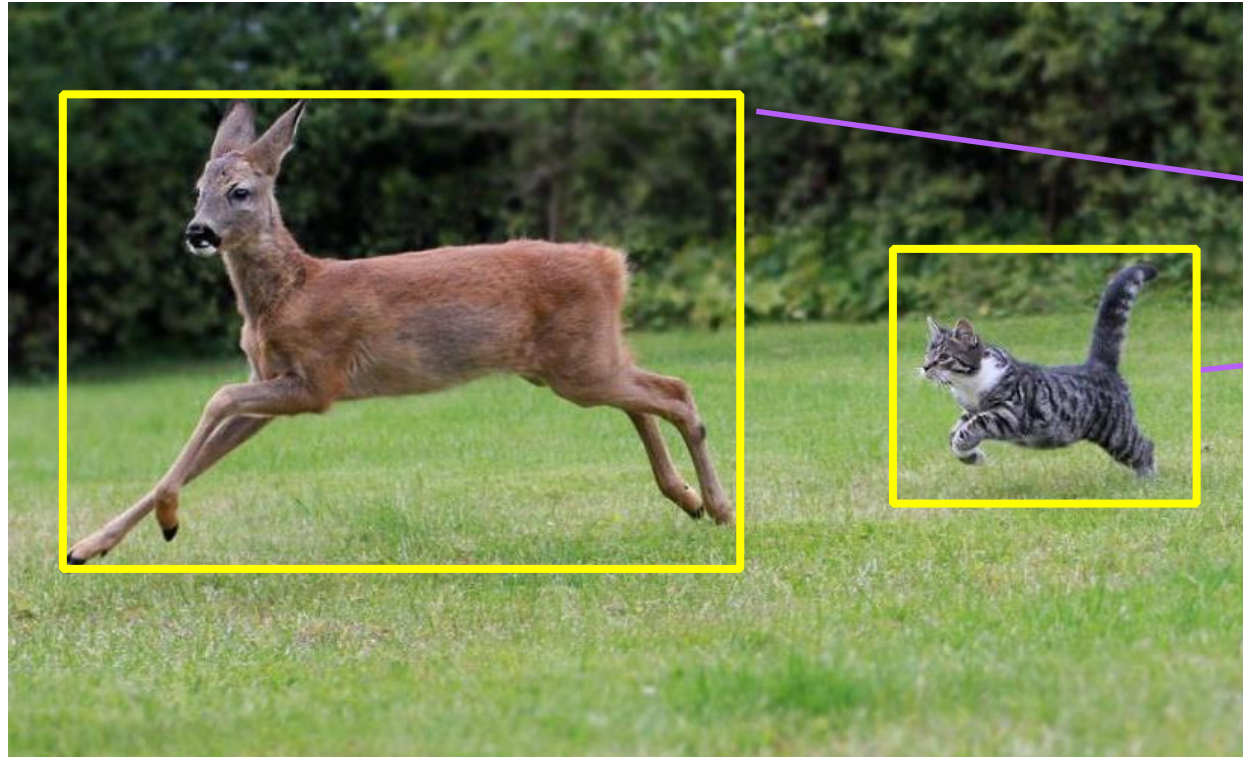
Image tagging



deer  
cat  
trees  
grass

# What is Visual Recognition?

## Object detection



deer

cat

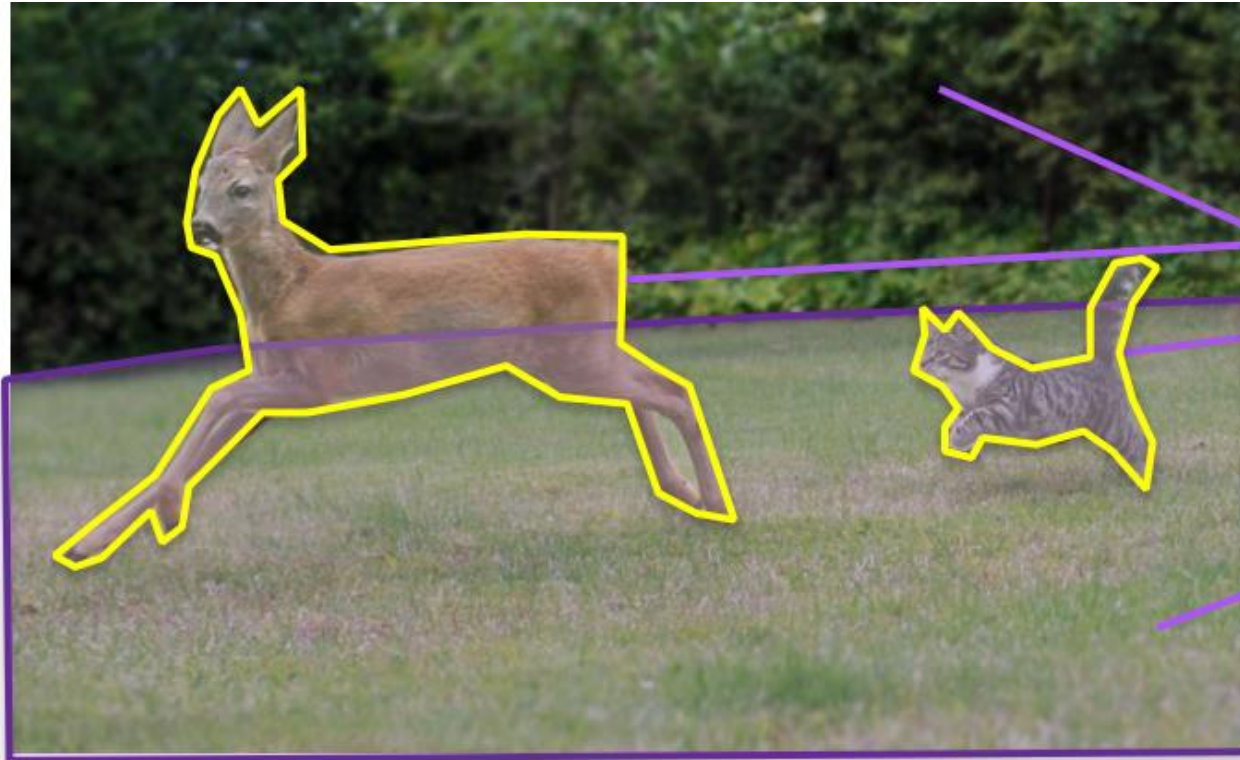
trees

grass



# What is Visual Recognition?

## Object segmentation



deer

cat

trees

grass

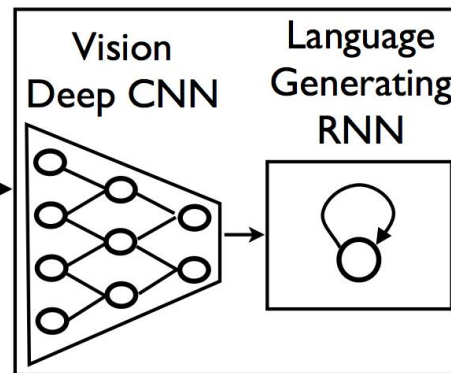
# Pushing the Limits of Visual Recognition

Reasoning about Language!



a cat is chasing a  
young deer on  
the grass

# Vision & Language – Visual Captioning



**A group of people shopping at an outdoor market.**

**There are many vegetables at the fruit stand.**

Vinyals et al. 2015



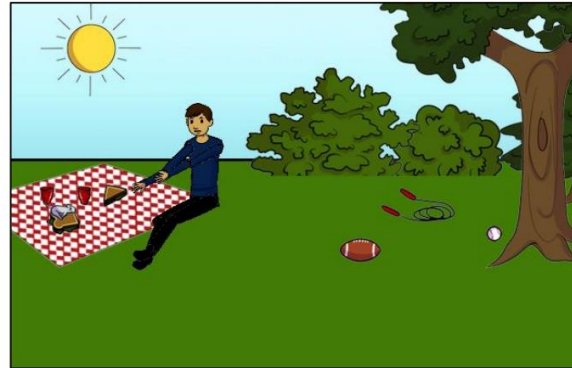
# Vision & Language – Visual QA



What color are her eyes?  
What is the mustache made of?



How many slices of pizza are there?  
Is this a vegetarian pizza?



Is this person expecting company?  
What is just under the tree?



Does it appear to be rainy?  
Does this person have 20/20 vision?

Agrawal et al. 2015



# Vision & Language – Textual Grounding

---

The two girls in hats in the middle



Rohrbach et al. 2016

# Vision & Language – Text to Image

---

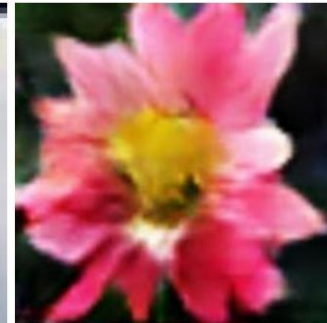
This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face



This bird is white with some black on its head and wings, and has a long orange beak



This flower has overlapping pink pointed petals surrounding a ring of short yellow filaments



Zhang et al. 2016

# Vision & Language



Vision-to-Language  
Generation



Visual Captioning



Image Captioning

Video Captioning

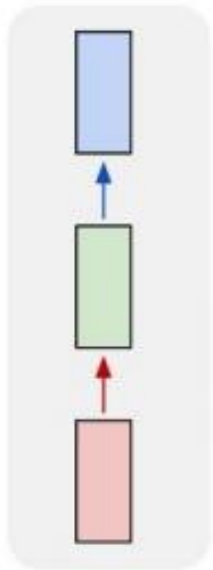
# Image Captioning

---

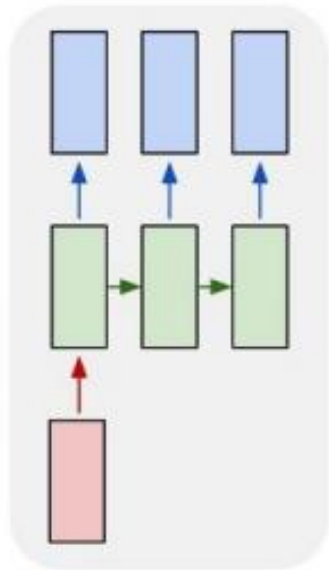


# Recall: Recurrent Neural Network

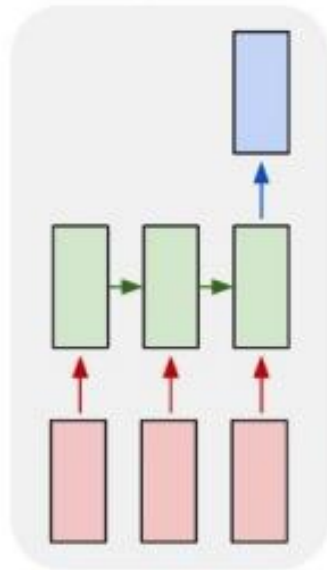
one to one



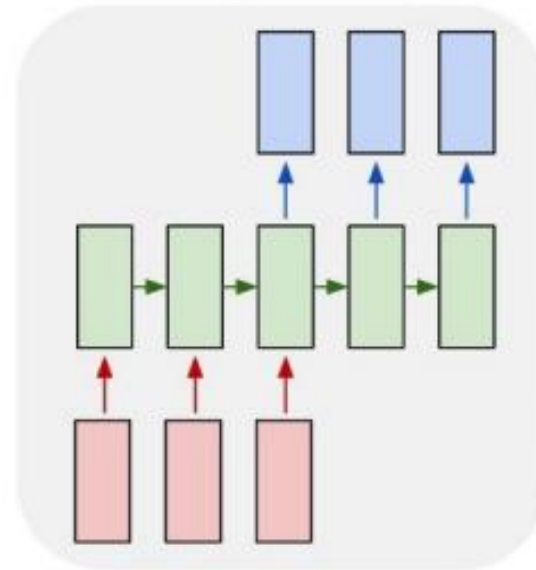
one to many



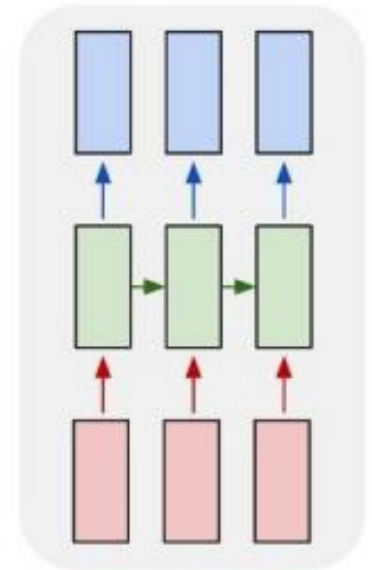
many to one



many to many



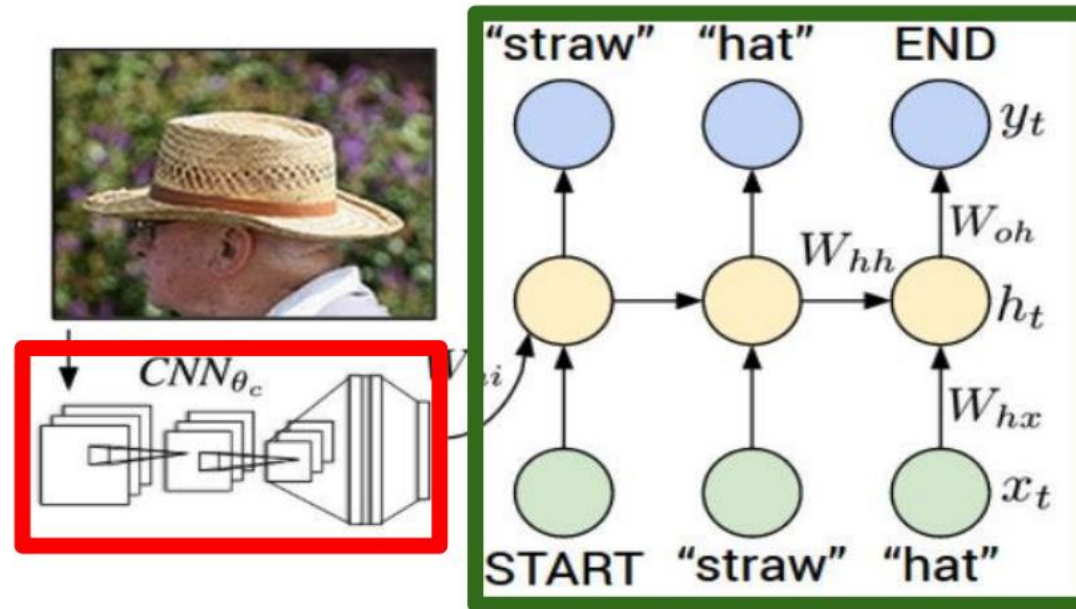
many to many



↖ e.g. **Image Captioning**  
image -> sequence of words

# Image Captioning

## Recurrent Neural Network



## Convolutional Neural Network



test image

[This image is CC0 public domain](#)

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

FC-1000

softmax



test image



image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

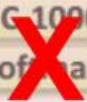
FC-1000

softmax



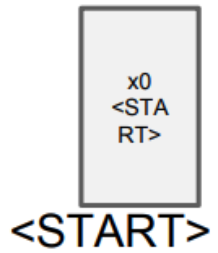
[This image is CC0 public domain](#)

test image





test image



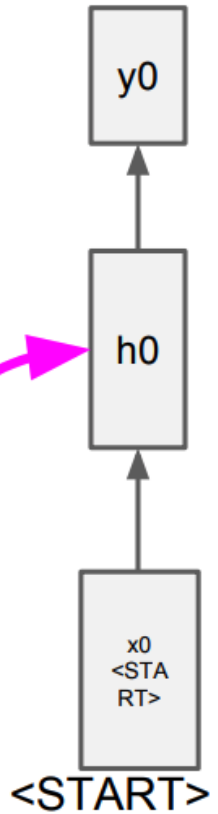


V



test image

Wih



before:

$$h = \tanh(W_{xh} * x + W_{hh} * h)$$

now:

$$h = \tanh(W_{xh} * x + W_{hh} * h + W_{ih} * v)$$

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

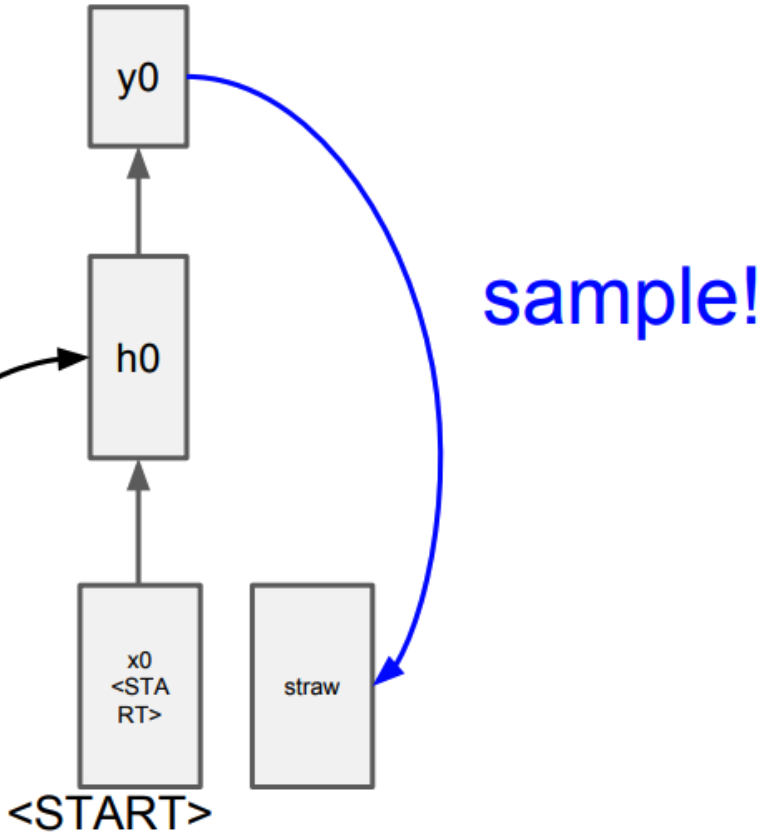
maxpool

FC-4096

FC-4096



test image





image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

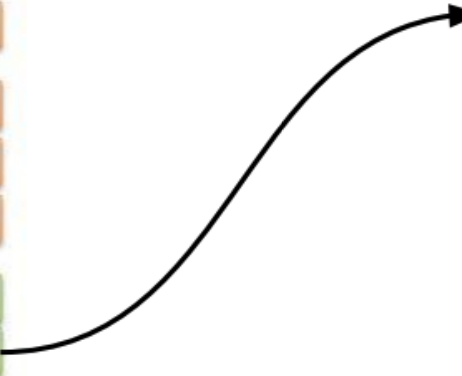
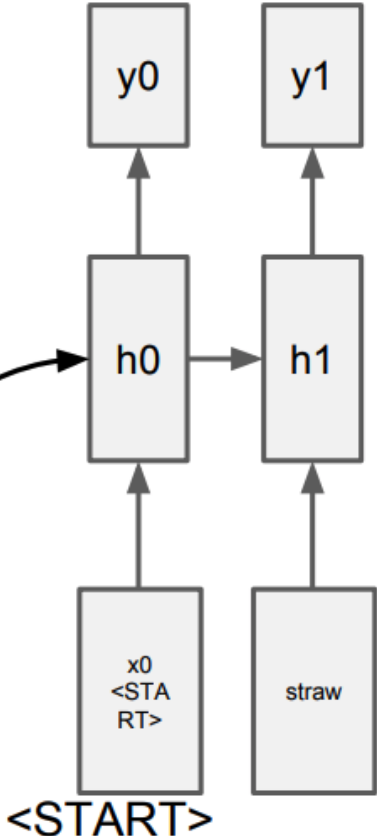
maxpool

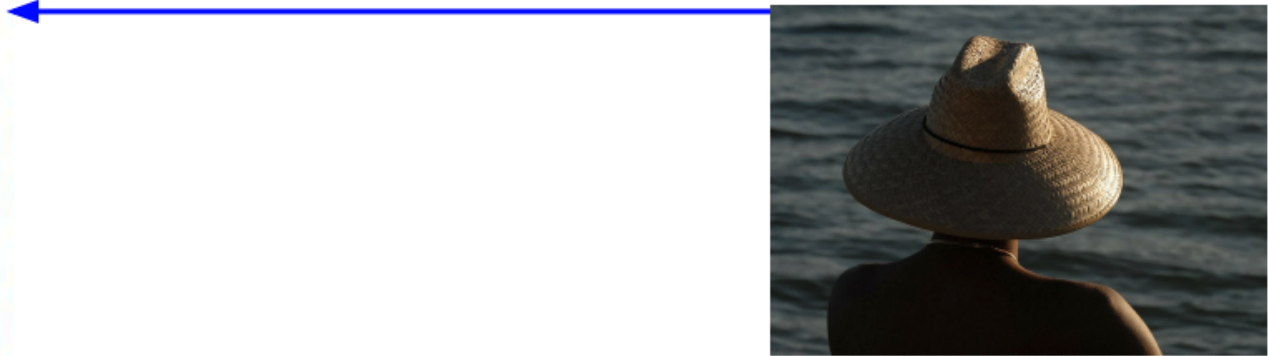
FC-4096

FC-4096

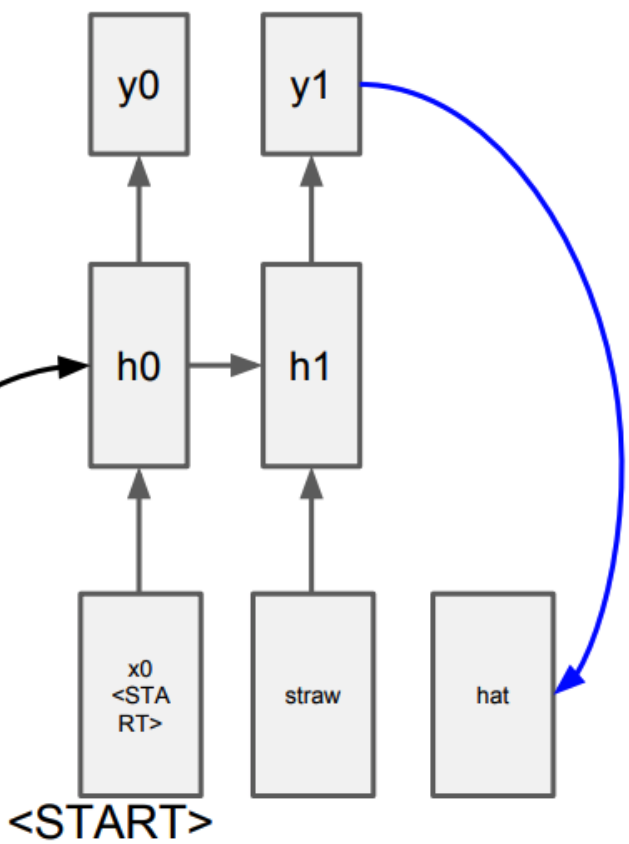


test image





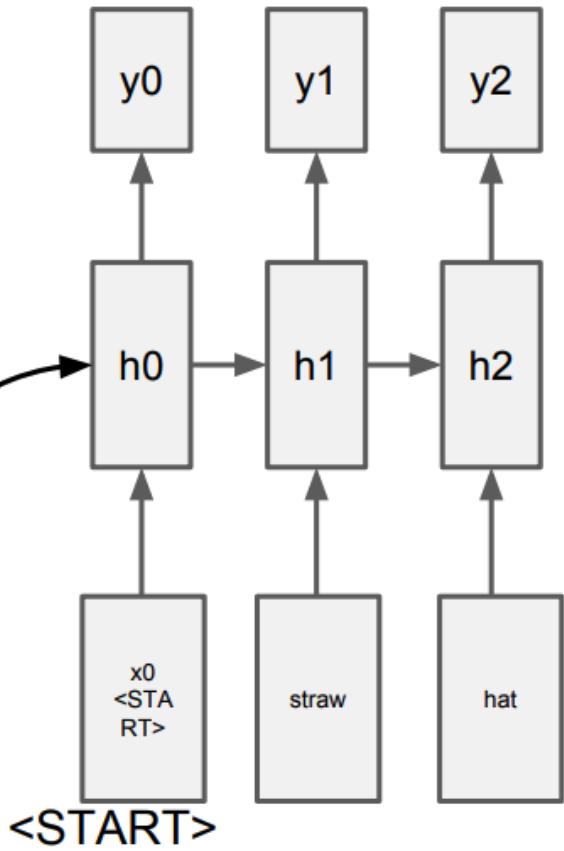
test image



sample!

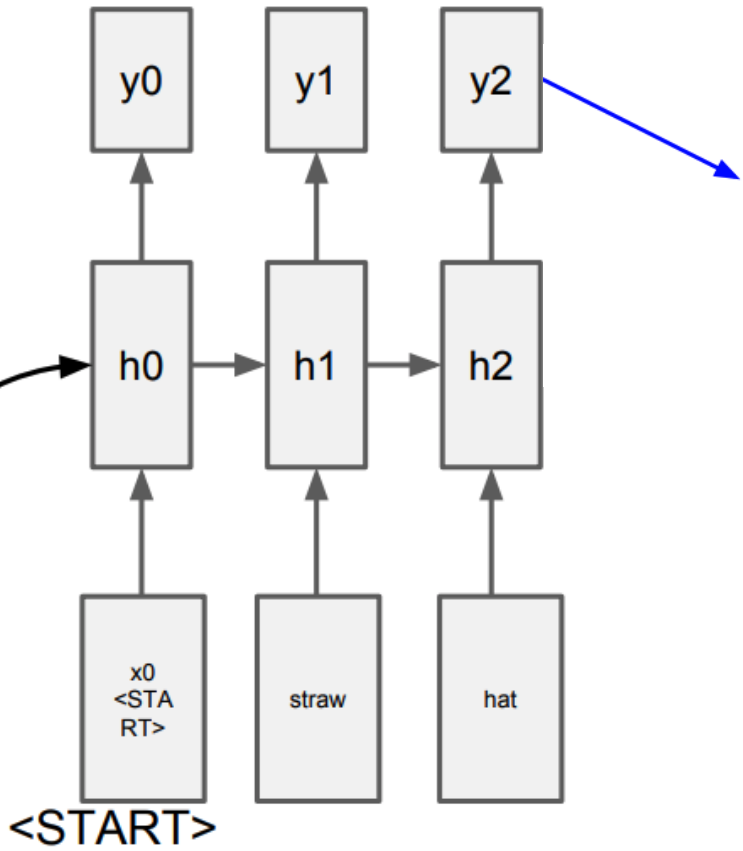


test image





test image



sample  
<END> token  
=> finish.



# Image Captioning: Example Results



*A cat sitting on a suitcase on the floor*



*A cat is sitting on a tree branch*



*A dog is running in the grass with a frisbee*



*A white teddy bear sitting in the grass*



*Two people walking on the beach with surfboards*



*A tennis player in action on the court*



*Two giraffes standing in a grassy field*



*A man riding a dirt bike on a dirt track*



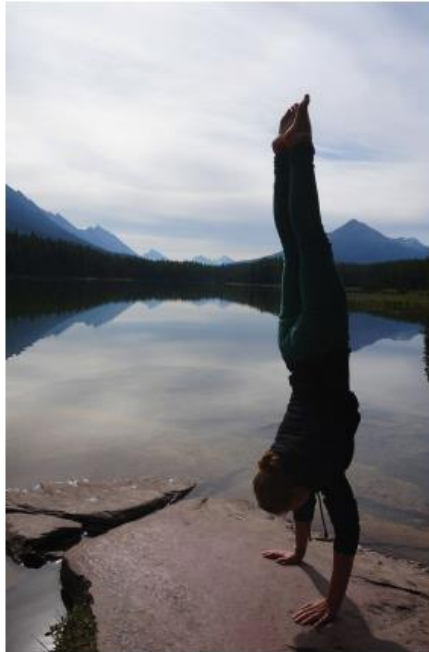
# Image Captioning: Failure Cases



*A woman is holding a cat in her hand*



*A person holding a computer mouse on a desk*



*A woman standing on a beach holding a surfboard*



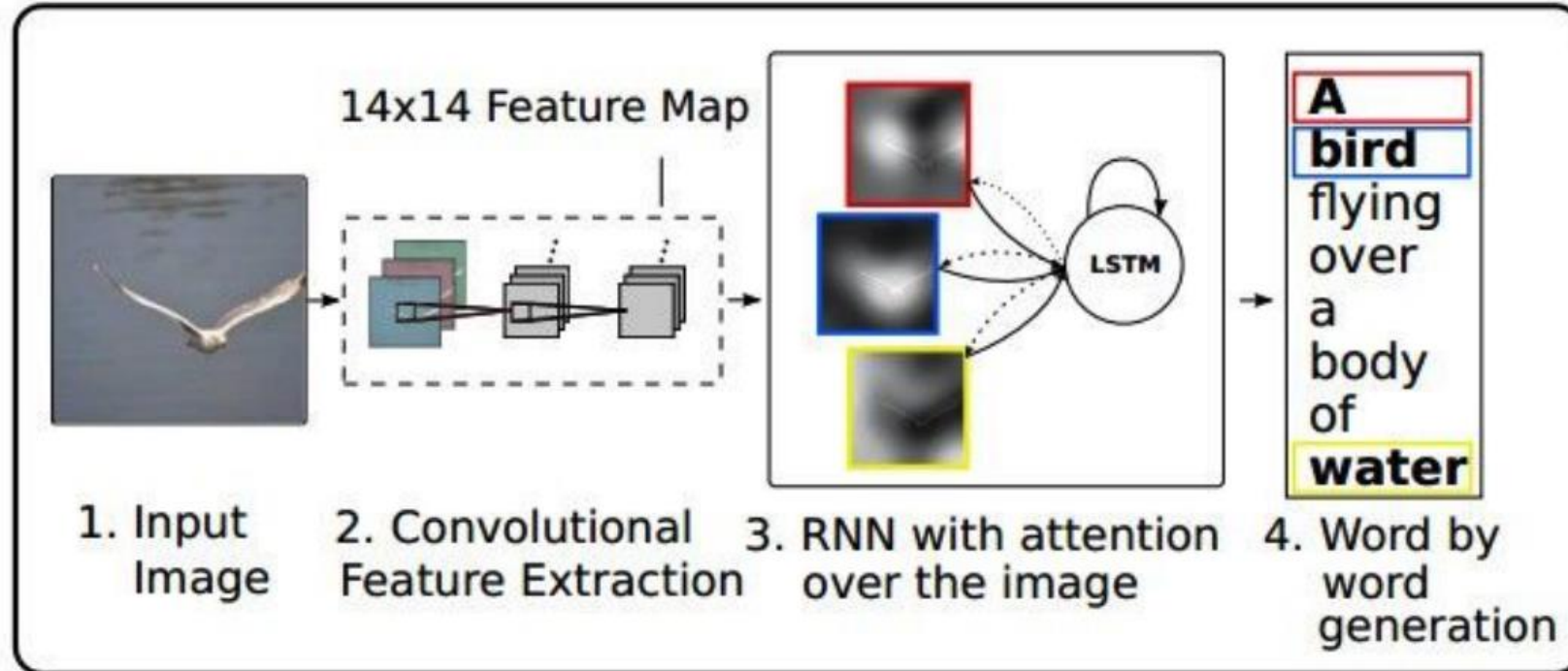
*A bird is perched on a tree branch*



*A man in a baseball uniform throwing a ball*

# Image Captioning with Attention

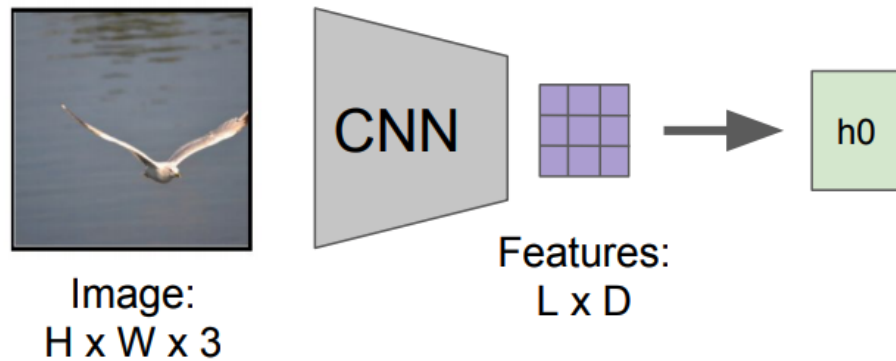
RNN focuses its attention at a different spatial location when generating each word



Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

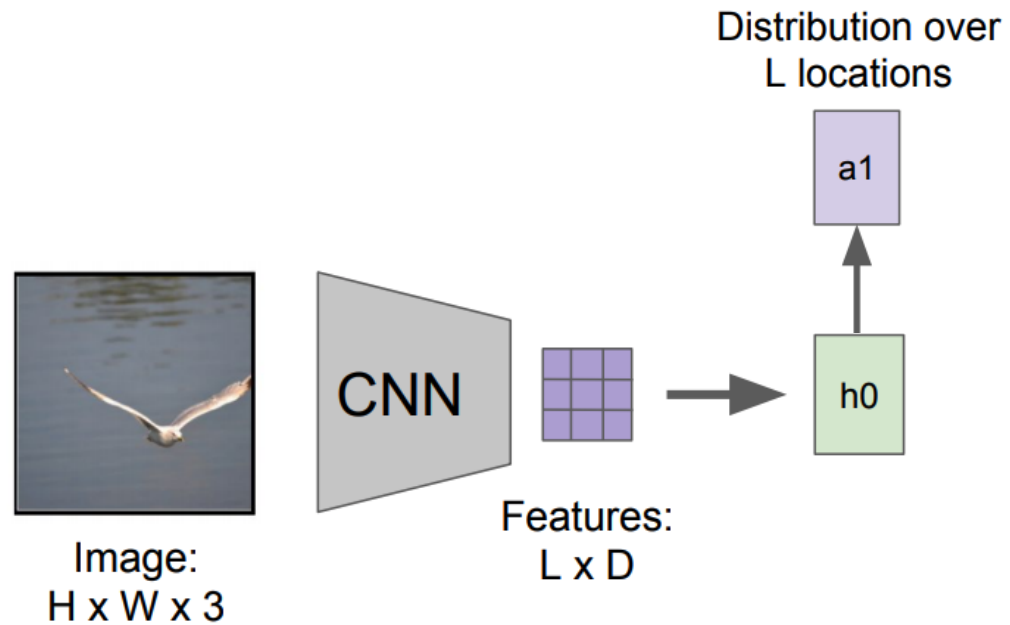
# Image Captioning with Attention

---



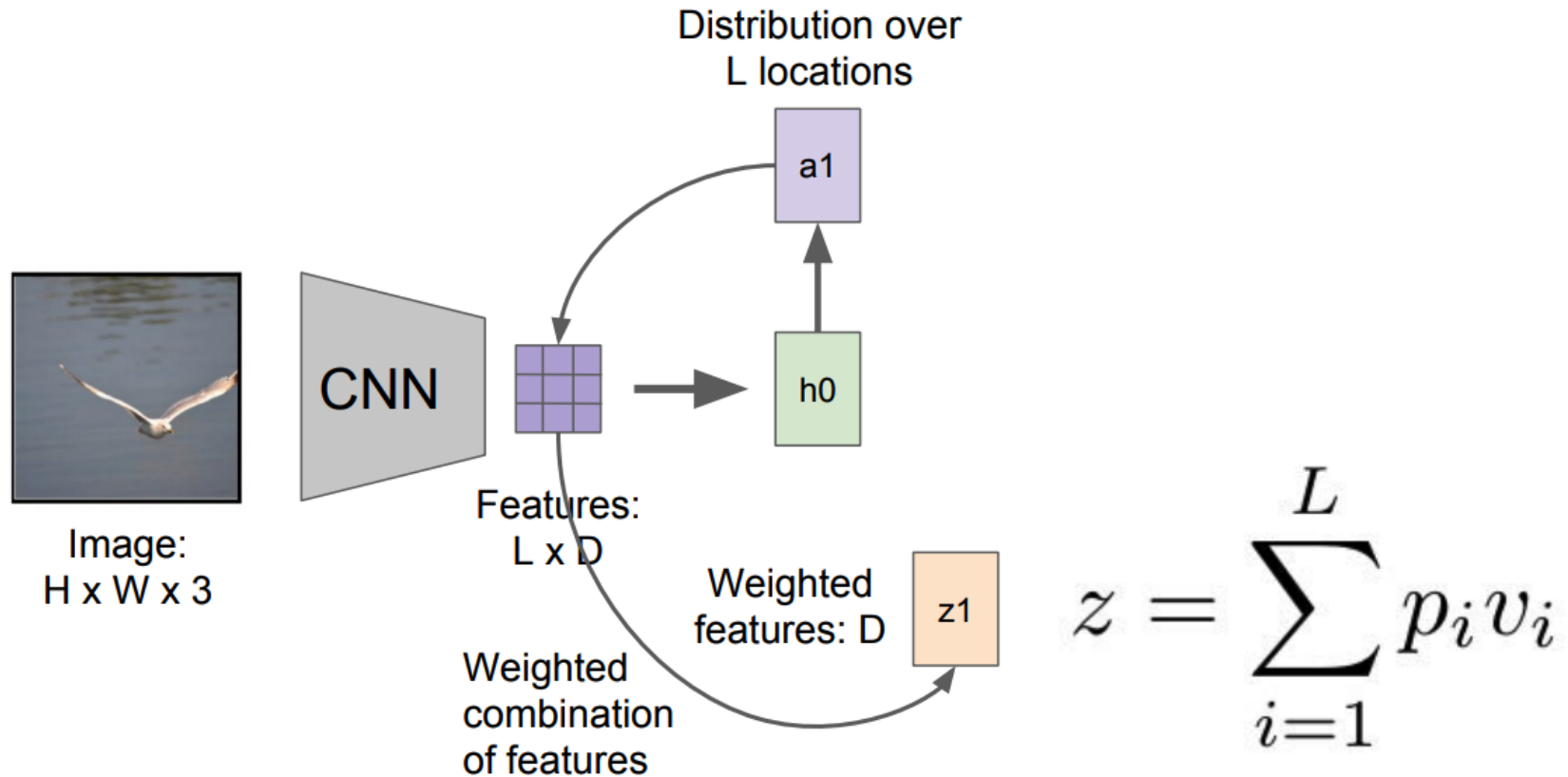
Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with Attention



Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

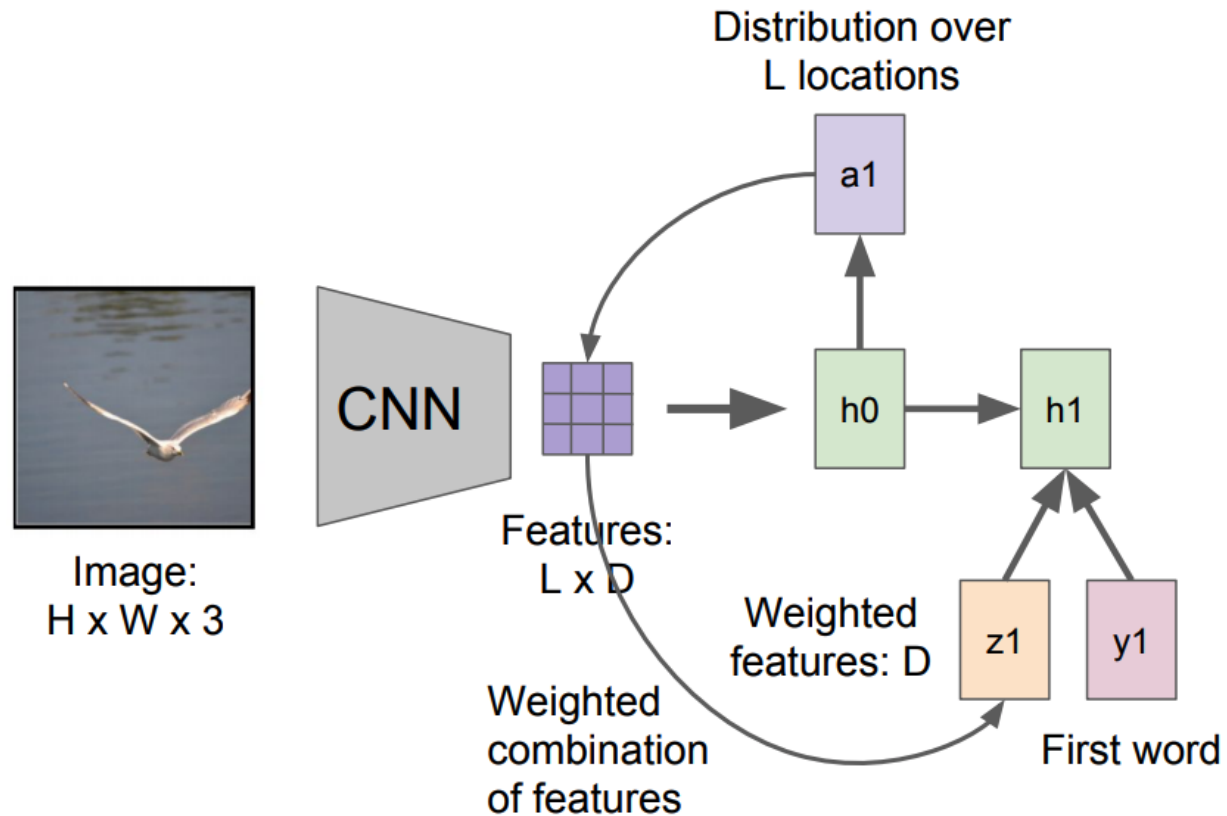
# Image Captioning with Attention



Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

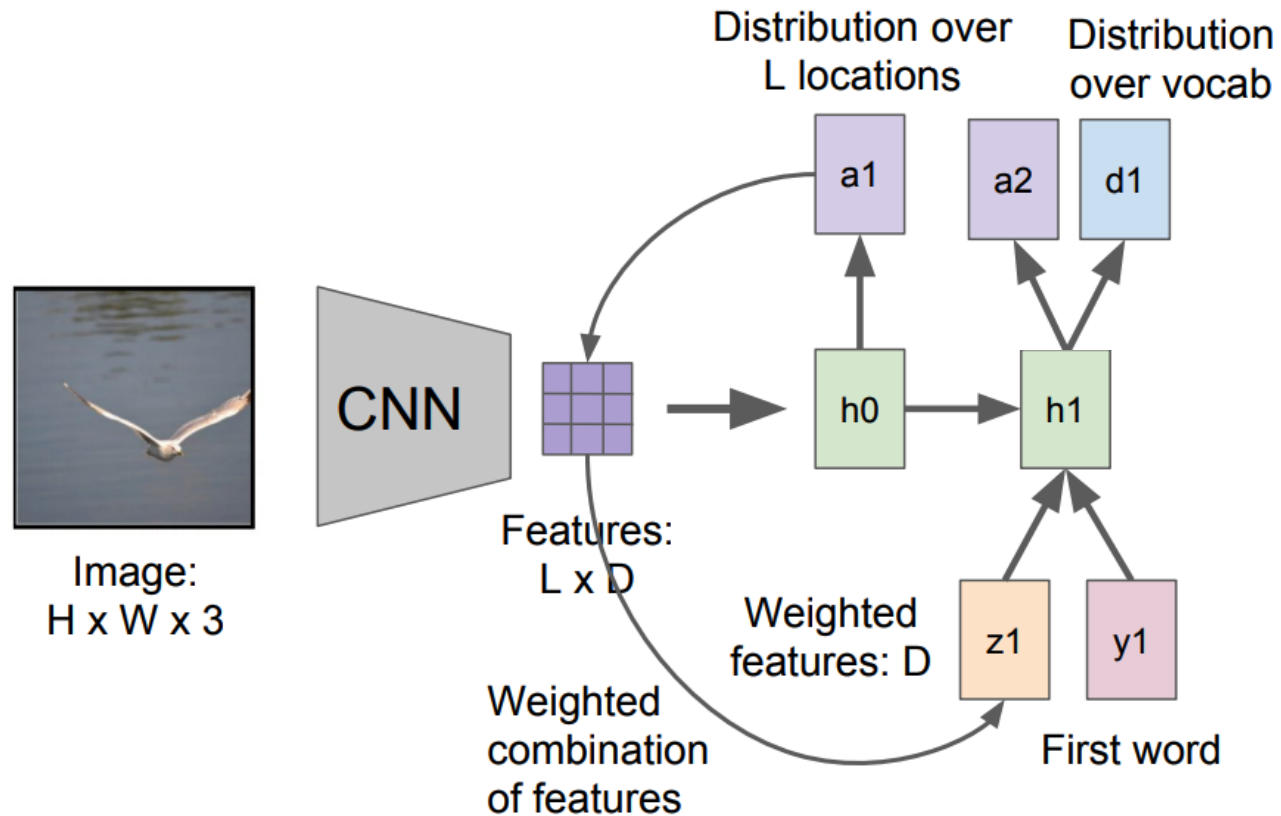


# Image Captioning with Attention



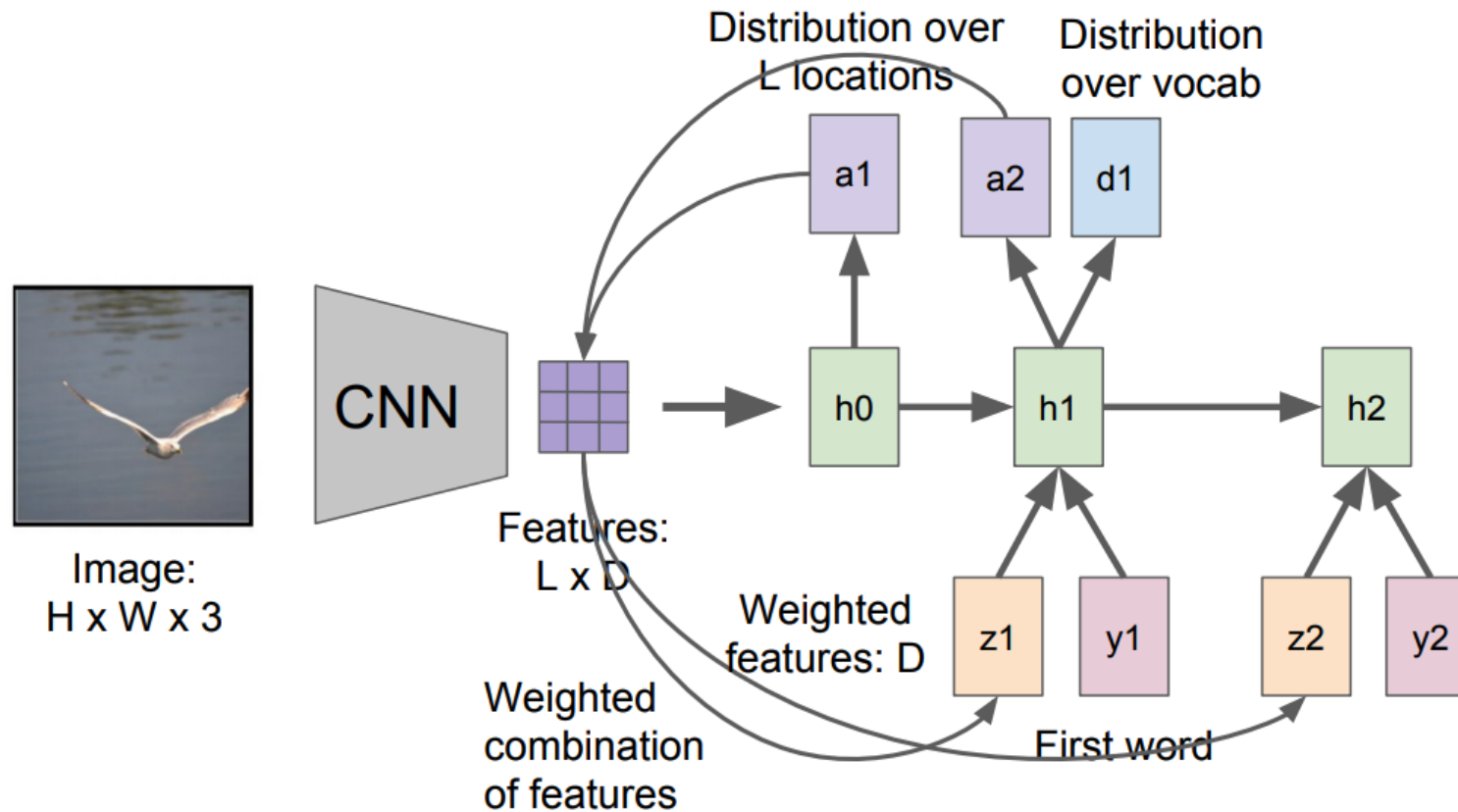
Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with Attention



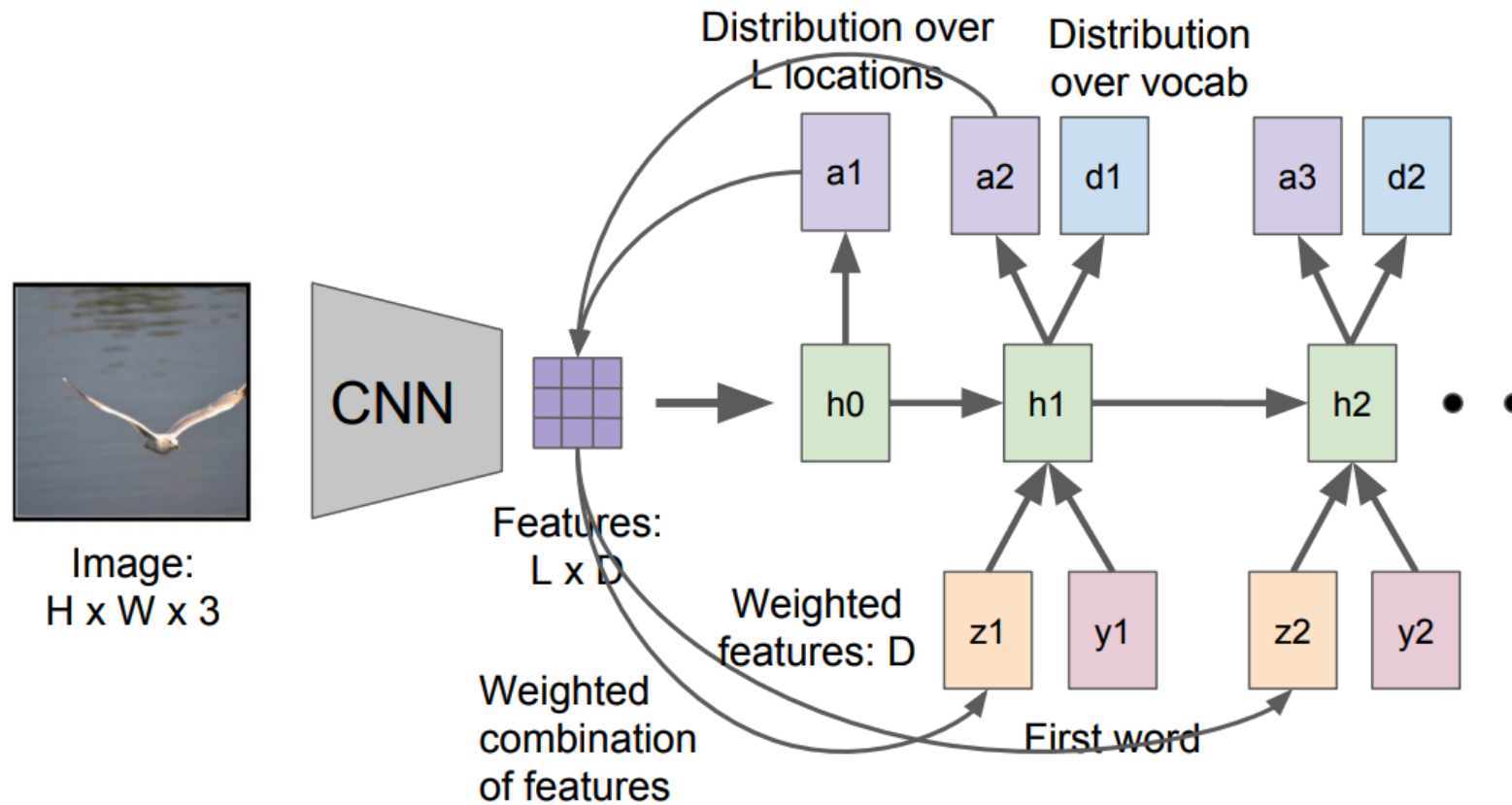
Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with Attention



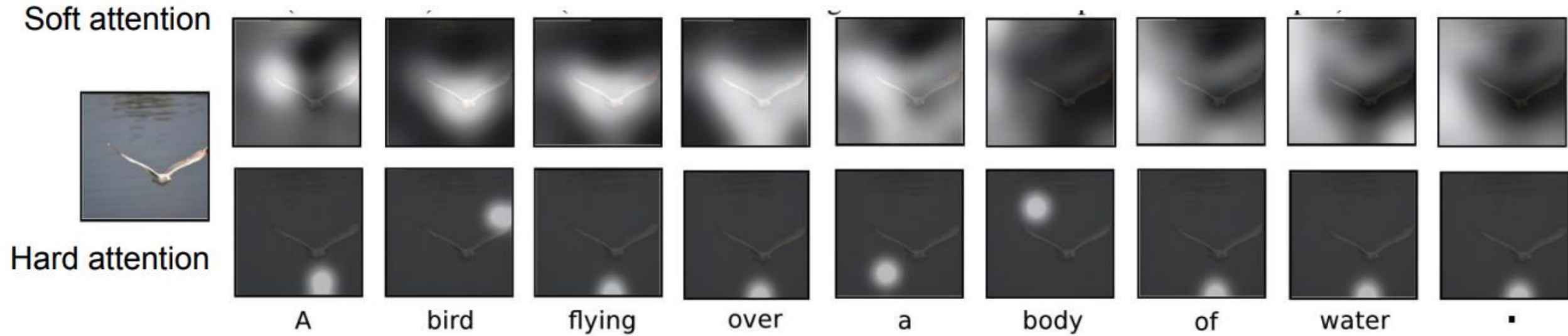
Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with Attention



Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with Attention



Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015



# Image Captioning with Attention



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



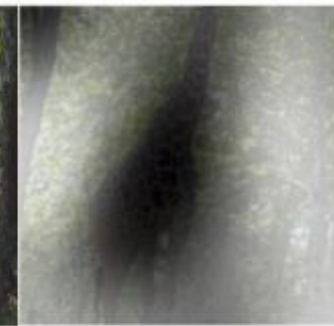
A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.



Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Video Captioning

---

# Video Captioning

Single-sentence  
Generation



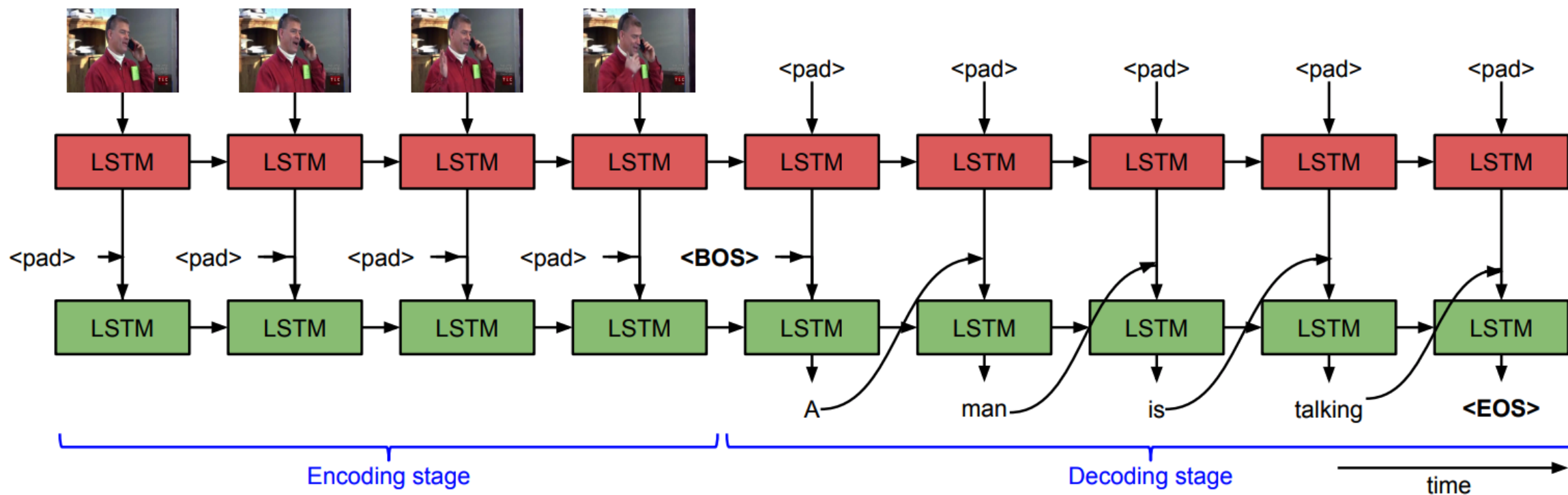
*A dog is playing in a bowl.*

Paragraph  
Generation



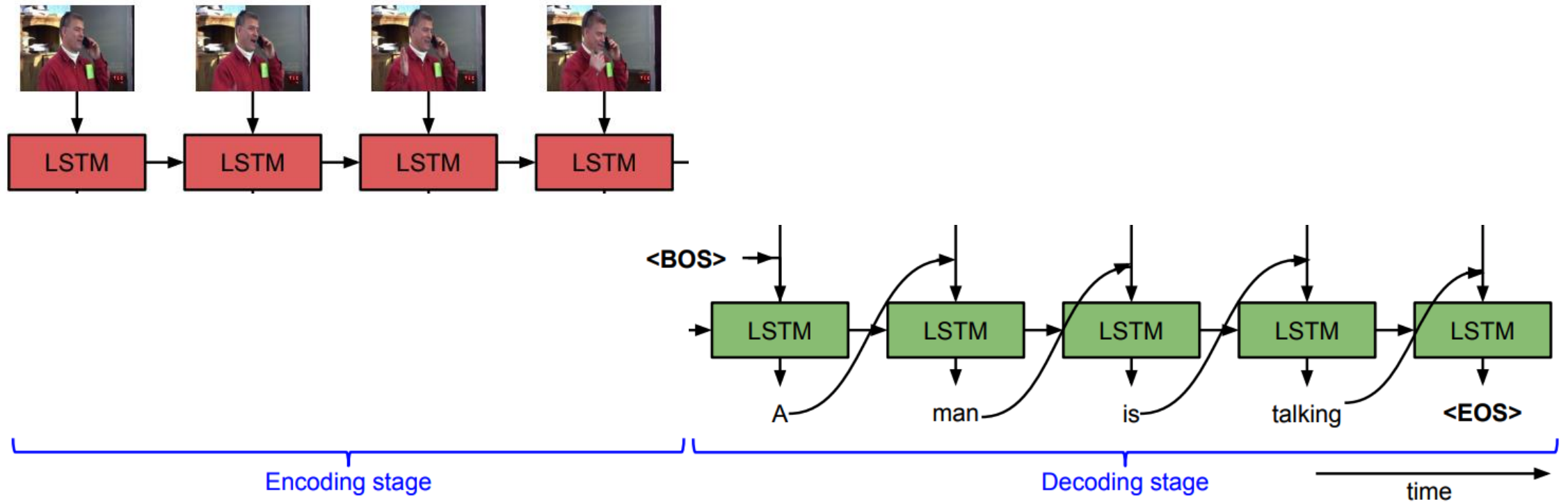
*The person peeled the fruit.  
The person put the fruit in the bowl.  
The person sliced the orange.  
The person put the pieces in the plate.  
The person rinsed the plate in the sink.*

# Sequence-to-sequence Model



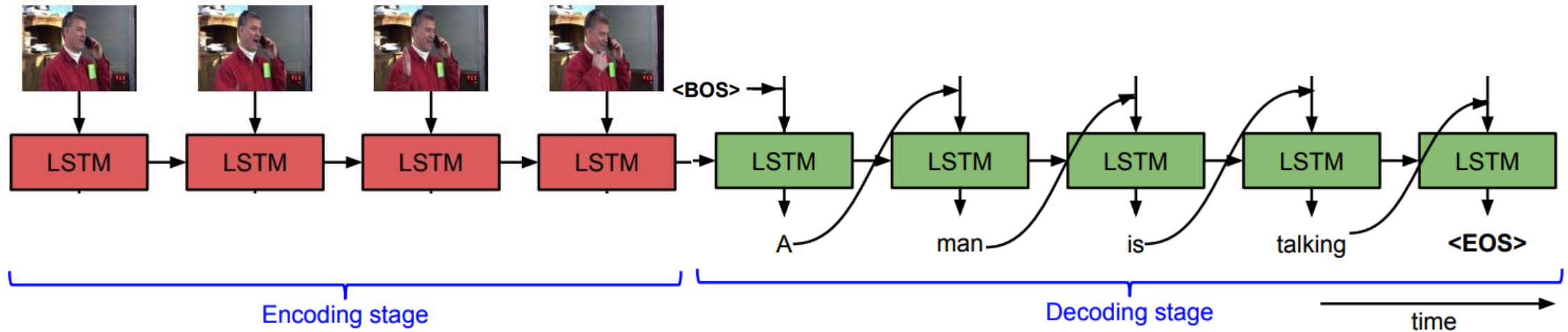
Venugopalan et al, "Sequence to Sequence – Video to Text", ICCV 2015

# Sequence-to-sequence Model





# Sequence-to-sequence Model



# Takeaway

---

- Vision & Language tasks require deeper understanding of the images/videos
- Transfer Learning
- Visual Captioning
  - CNN-RNN architecture
  - Sequence-to-sequence models
  - Attention Mechanism

# Thank You !

Xin Wang

[xwang@cs.ucsb.edu](mailto:xwang@cs.ucsb.edu)

<http://www.cs.ucsb.edu/~xwang>